

# Linear Models

EECS 442 – Jeong Joon Park  
Winter 2024, University of Michigan

# Roadmap

- Basics (Done)
- Image Stitching/Warping (Done)
- Learning-based Vision + Generative Models  
(through end of March)
- 3D Vision (April)
- Applications (Last Week)

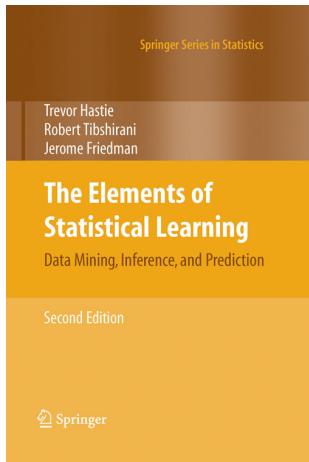
# Next Few Classes

- Machine Learning (ML) Crash Course
- I can't cover everything
- If you can, take a ML course or *learn online*
- ML really won't solve all problems and is incredibly dangerous if misused
- But ML is a powerful tool and not going away

# Terminology

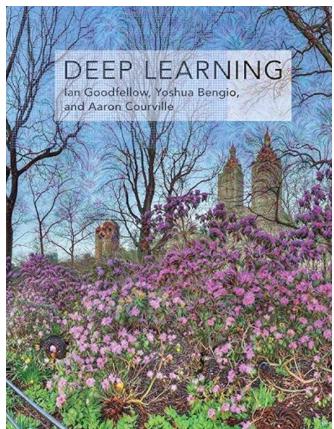
- ML is incredibly messy terminology-wise.
- Most things have lots of names.
- I will try to write down multiple of them so if you see it later you'll know what it is.

# Pointers



Useful book (Free too!):  
The Elements of Statistical Learning  
Hastie, Tibshirani, Friedman

<https://web.stanford.edu/~hastie/ElemStatLearn/>



Useful set of data:  
UCI ML Repository  
<https://archive.ics.uci.edu/ml/datasets.html>

Deep Learning  
I. GoodFellow et al.  
<https://www.deeplearningbook.org/>

# Machine Learning (ML)

- Goal: make “sense” of data
- Overly simplified version: transform vector  $\mathbf{x}$  into vector  $\mathbf{y}=T(\mathbf{x})$  that’s somehow better
- Potentially you fit  $T$  using pairs of datapoints and desired outputs  $(\mathbf{x}_i, \mathbf{y}_i)$ , or just using a set of datapoints  $(\mathbf{x}_i)$
- Always are trying to find some transformation that minimizes or maximizes some **objective function** or goal.

# Machine Learning

Input:  $x$

**Feature vector/Data point:**

Vector representation of datapoint. Each dimension or “**feature**” represents some aspect of the data.

Output:  $y$

**Label / target:**

Fixed length vector of desired output. Each dimension represents some aspect of the output data

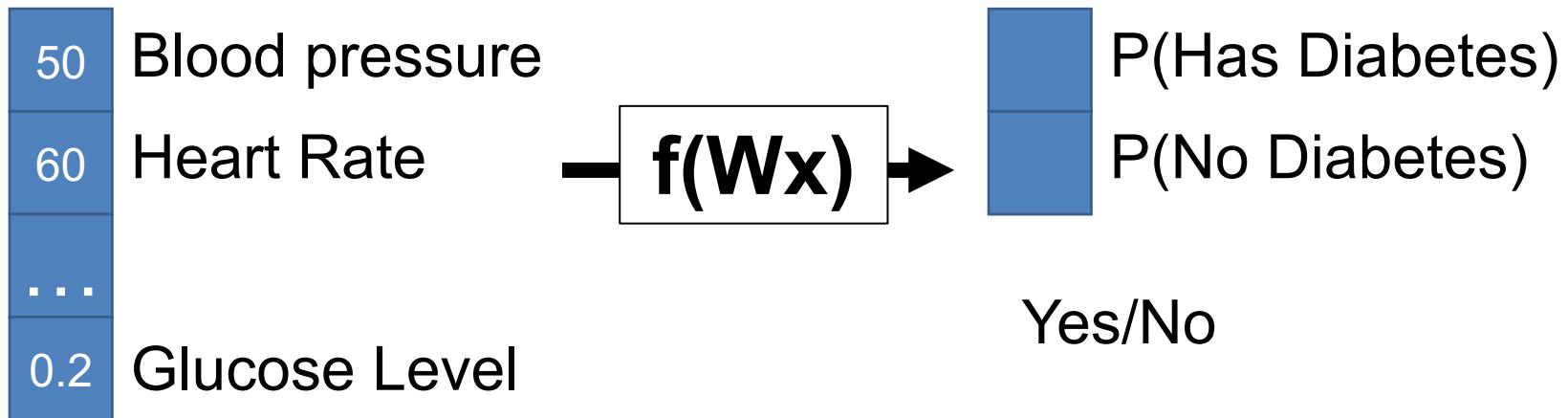
**Supervised:** we are given  $y$ .

**Unsupervised:** we are not, and make our own  $y$ s.

# Example – Health

Input:  $x$  in  $\mathbb{R}^N$

Output:  $y$

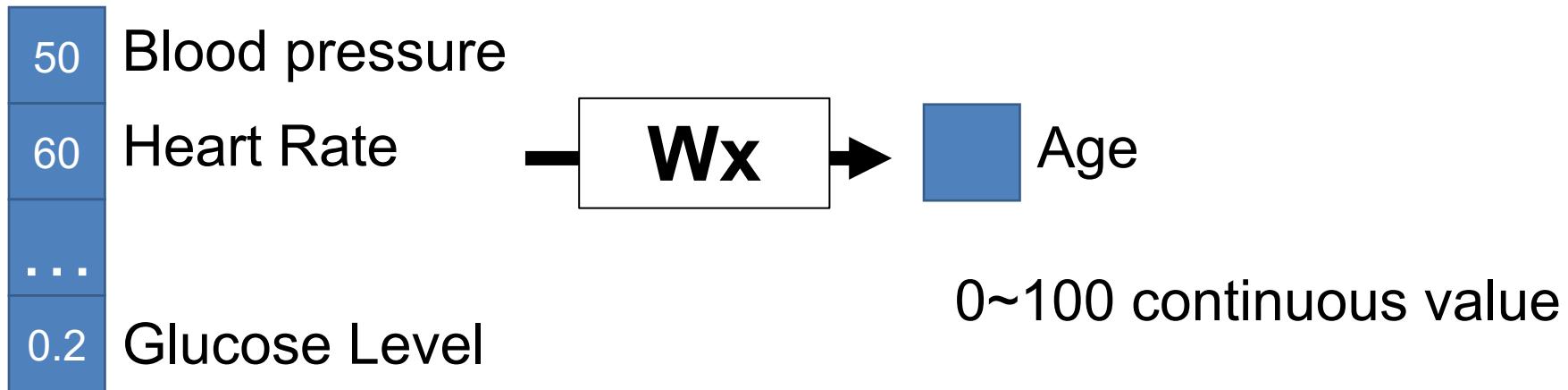


*Intuitive objective function:* Want correct category to be likely with our model.

# Example – Health

Input:  $x$  in  $\mathbb{R}^N$

Output:  $y$

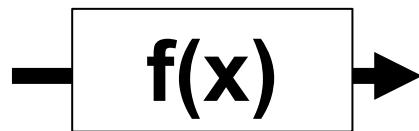


*Intuitive objective function:* Want our prediction of age to be “close” to true age.

# Example – Health

Input:  $x$  in  $\mathbb{R}^N$

50	Blood pressure
60	Heart Rate
...	
0.2	Glucose Level



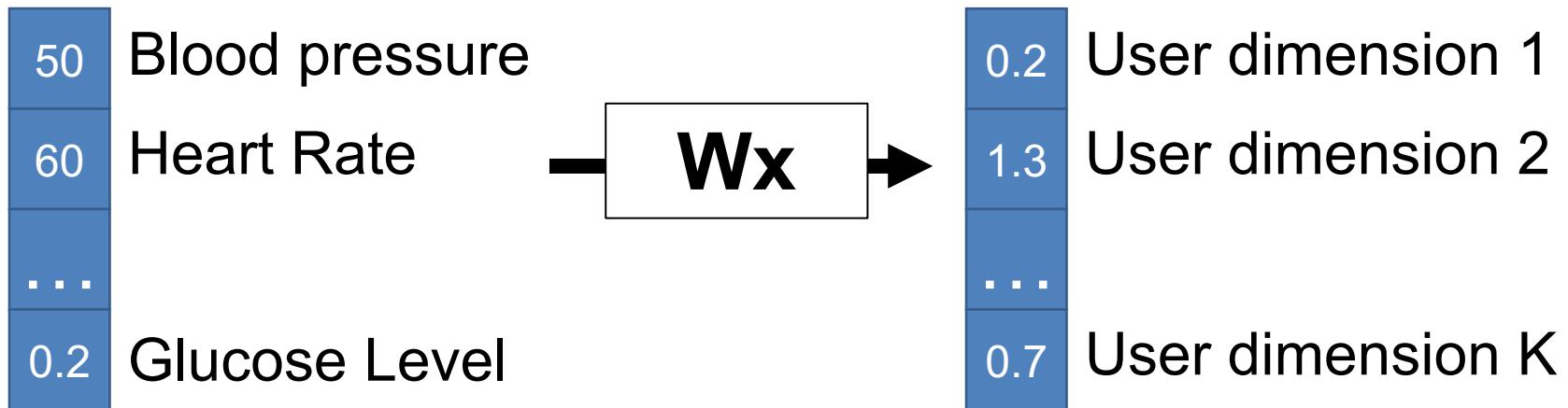
Output: **discrete  $y$**   
**(unsupervised)**

*Intuitive objective function:* Want to find  $K$  groups that explain the data we see.

# Example – Health

Input:  $x$  in  $\mathbb{R}^N$

Output: **continuous  $y$**   
**(discovered)**

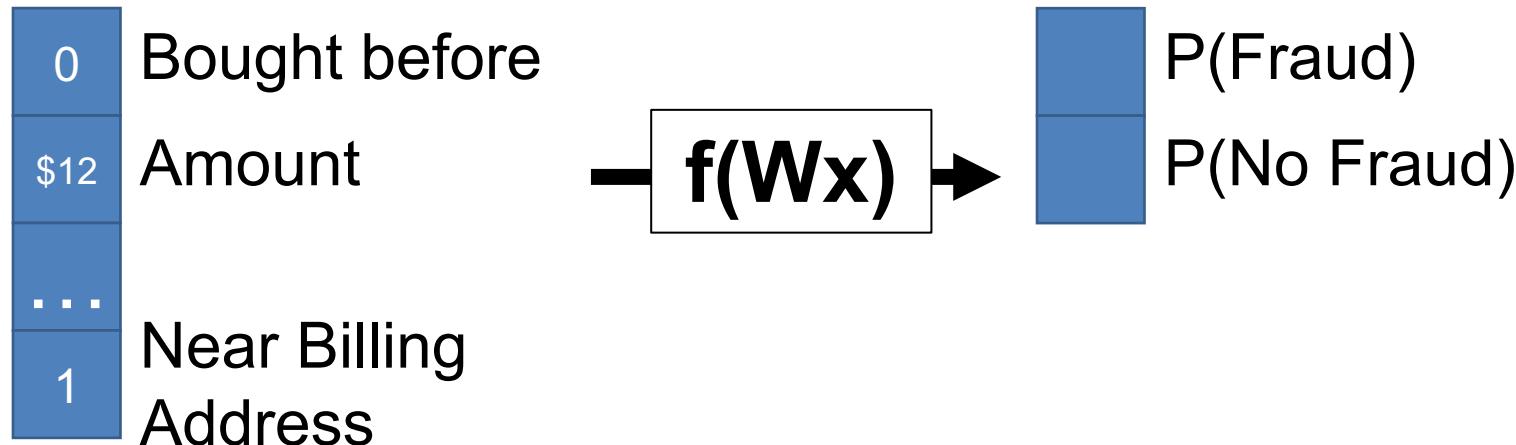


*Intuitive objective function:* Want to  $K$  dimensions (often two) that are easier to understand but capture the variance of the data.

# Example – Credit Card Fraud

Input:  $x$  in  $\mathbb{R}^N$

Output:  $y$

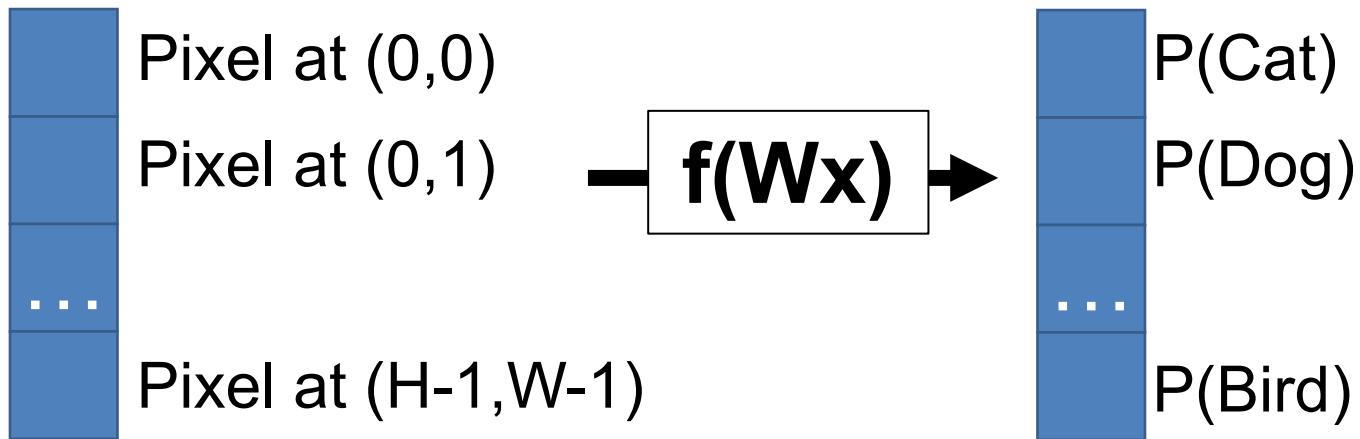


*Intuitive objective function:* Want correct category to be likely with our model.

# Example – Computer Vision

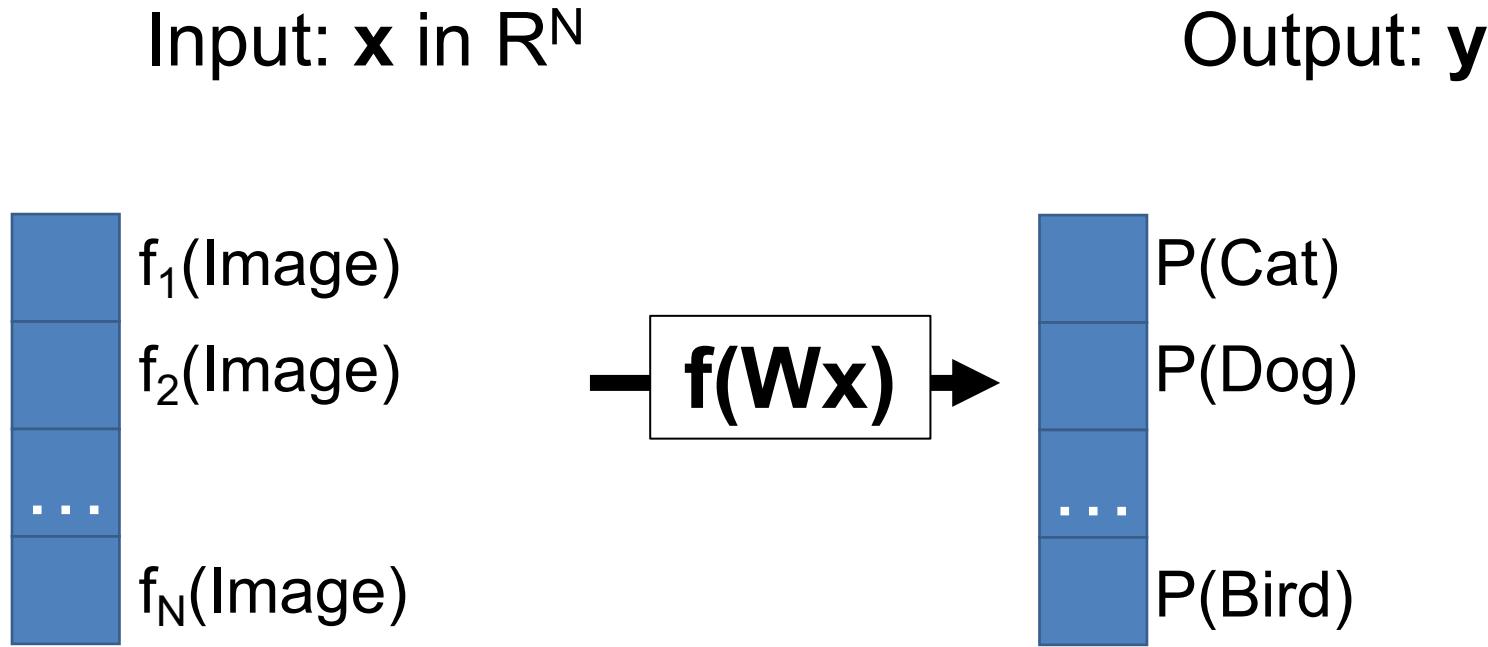
Input:  $\mathbf{x}$  in  $\mathbb{R}^N$

Output:  $\mathbf{y}$



*Intuitive objective function:* Want correct category to be likely with our model.

# Example – Computer Vision



*Intuitive objective function:* Want correct category to be likely with our model.

# Abstractions

- Throughout, assume we've converted data into a fixed-length feature vector. There are well-designed ways for doing this.
- But remember it could be big!
  - Image (e.g., 224x224x3): 151K dimensions
  - Patch (e.g., 32x32x3) in image: 3072 dimensions

# ML Problems in Vision



# ML Problem Examples in Vision

**Supervised  
(Data+Labels)**

**Unsupervised  
(Just Data)**

**Discrete  
Output**

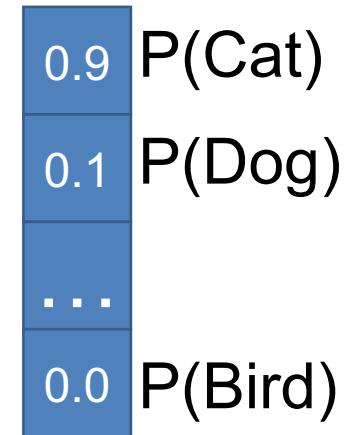
**Classification/  
Categorization**

**Continuous  
Output**

# ML Problem Examples in Vision

## Categorization/Classification

Binning into K mutually-exclusive categories



# ML Problem Examples in Vision

**Supervised  
(Data+Labels)**

**Unsupervised  
(Just Data)**

**Discrete  
Output**

Classification/  
Categorization

**Continuous  
Output**

**Regression**

# ML Problem Examples in Vision

## Regression

Estimating continuous variable(s)



3.6  
kg

Cat weight

# ML Problem Examples in Vision

**Supervised  
(Data+Labels)**

**Unsupervised  
(Just Data)**

**Discrete  
Output**

Classification/  
Categorization

**Clustering**

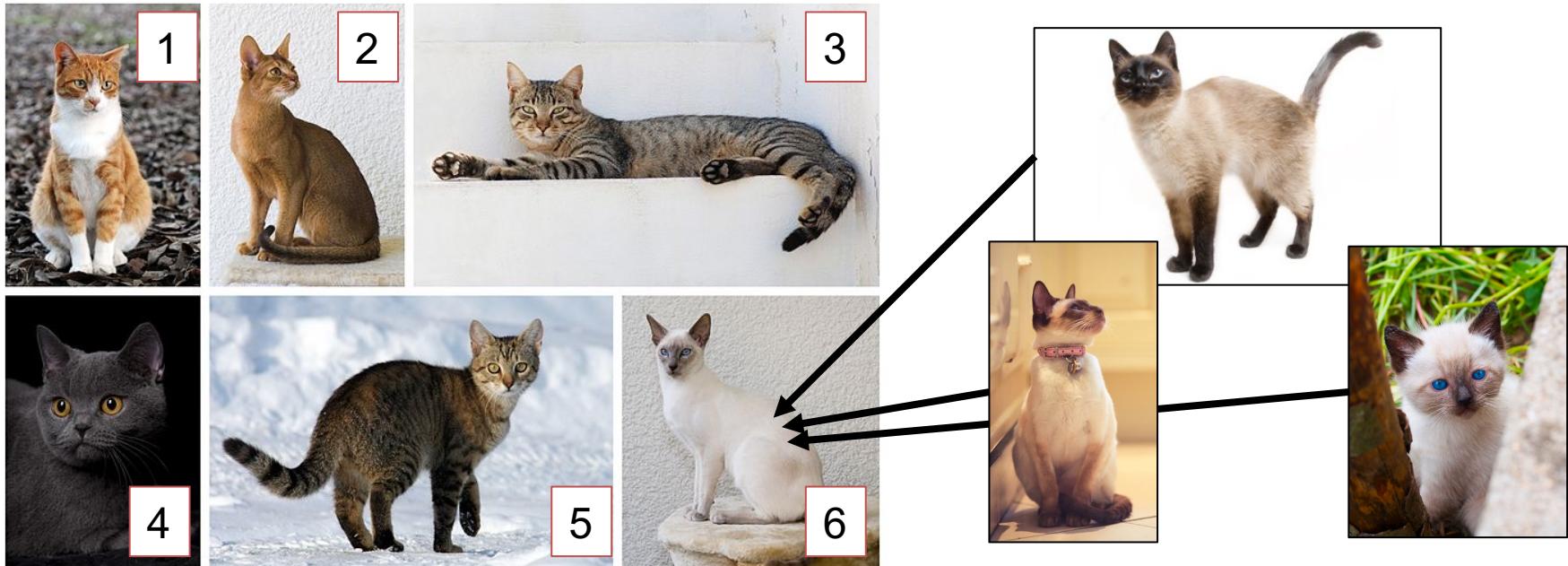
**Continuous  
Output**

Regression

# ML Problem Examples in Vision

## Clustering

Given a set of cats, automatically discover clusters or categories.



# ML Problem Examples in Vision

**Supervised  
(Data+Labels)**

**Unsupervised  
(Just Data)**

**Discrete  
Output**

Classification/  
Categorization

Clustering

**Continuous  
Output**

Regression

**Dimensionality  
Reduction**

# ML Problem Examples in Vision

## Dimensionality Reduction

Find dimensions that best explain  
the whole image/input



Cat size in  
image

Location of  
cat in image

For ordinary images, this is currently a challenging task. For certain images (e.g., faces, this works reasonably well)

# Practical Example

- ML has a tendency to be mysterious
- Let's start with least-squares
- One thing to remember:
  - $N \text{ eqns}, < N \text{ vars} = \text{overdetermined (will have errors)}$
  - $N \text{ eqns}, N \text{ vars} = \text{exact solution}$
  - $N \text{ eqns}, > N \text{ vars} = \text{underdetermined (infinite solns)}$

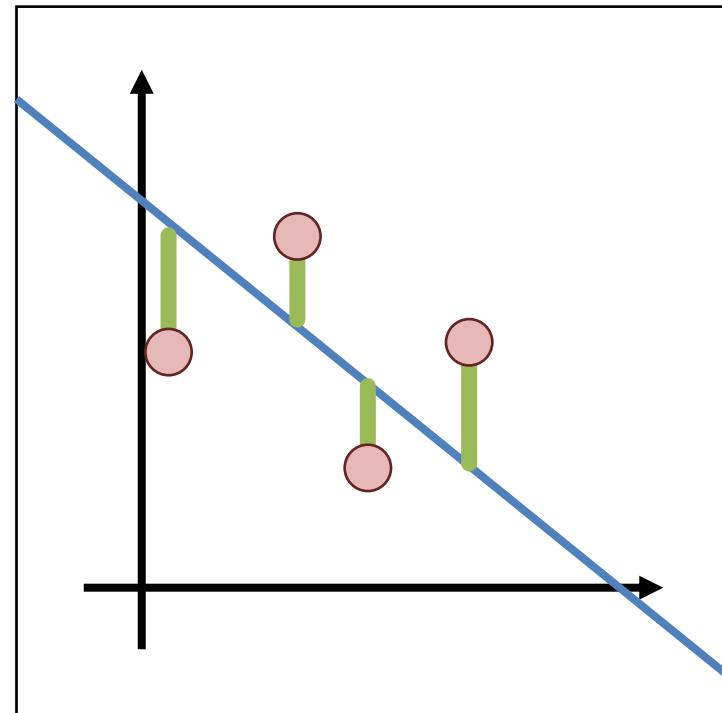
# Example – Least Squares

Let's make the world's **worst** weather model

Data:  $(x_1, y_1), (x_2, y_2), \dots, (x_k, y_k)$

Model:  $(m, b) \quad y_i = mx_i + b$   
Or  $(w) \quad y_i = w^T x_i$

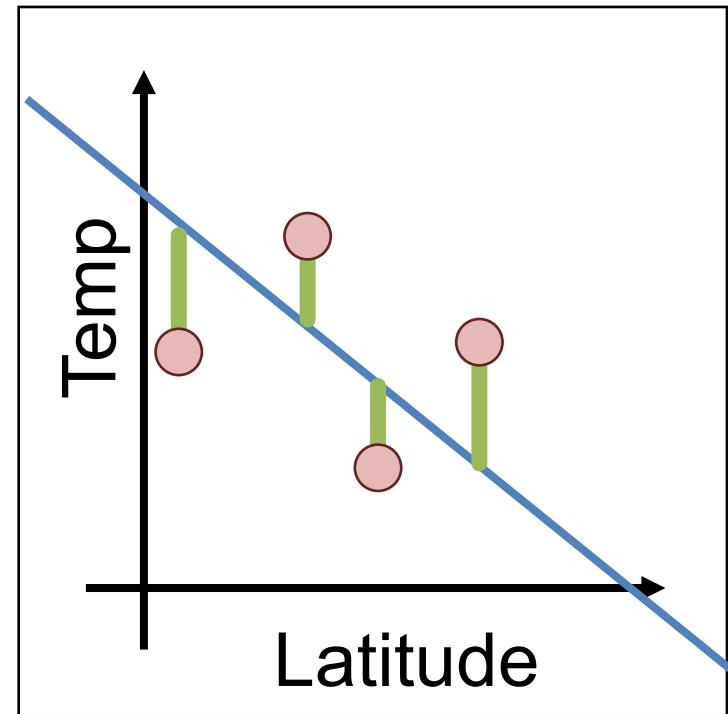
Objective function:  
 $(y_i - w^T x_i)^2$



# World's Worst Weather Model

Given latitude (distance above equator), predict temperature by fitting a line

<u>City</u>	<u>Latitude (°)</u>	<u>Temp (F)</u>
Ann Arbor	42	33
Washington, DC	39	38
Austin, TX	30	62
Mexico City	19	67
Panama City	9	83



# Example – Least Squares

$$\sum_{i=1}^k (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \rightarrow \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

**Output:**

Temperature

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_k \end{bmatrix}$$

**Inputs:**

Latitude, 1

$$\mathbf{X} = \begin{bmatrix} x_1 & 1 \\ \vdots & \vdots \\ x_k & 1 \end{bmatrix}$$

**Model/Weights:**

Latitude, “Bias”

$$\mathbf{w} = \begin{bmatrix} m \\ b \end{bmatrix}$$

# Example – Least Squares

$$\sum_{i=1}^k (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \rightarrow \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

**Output:**

Temperature

$$\mathbf{y} = \begin{bmatrix} 33 \\ \vdots \\ 83 \end{bmatrix}$$

**Inputs:**

Latitude, 1

$$\mathbf{X} = \begin{bmatrix} 42 & 1 \\ \vdots & \vdots \\ 9 & 1 \end{bmatrix}$$

**Model/Weights:**

Latitude, “Bias”

$$\mathbf{w} = \begin{bmatrix} m \\ b \end{bmatrix}$$

**Why do we add  
a one to the inputs?**

# Example – Least Squares

Training ( $\mathbf{x}_i, y_i$ ):

$$\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{or}$$

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2$$

**Loss function/objective:** evaluates correctness.  
Here: Squared L2 norm / Sum of Squared Errors

**Training/Learning/Fitting:** try to find model that  
*optimizes/minimizes* an objective / loss function

Optimal  $\mathbf{w}^*$  is

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

# Example – Least Squares

Training ( $\mathbf{x}_i, y_i$ ):

$$\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{or}$$

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2$$

Inference ( $\mathbf{x}$ ):

$$\mathbf{w}^T \mathbf{x} = w_1 x_1 + \cdots + w_F x_F$$

In our case:  $= mx + b$

**Testing/Inference:** Given a new output,  
what's the prediction?

# Least Squares: Learning

Data

<u>City</u>	<u>Latitude</u>	<u>Temp</u>
Ann Arbor	42	33
Washington, DC	39	38
Austin, TX	30	62
Mexico City	19	67
Panama City	9	83

Model

$$\text{Temp} = -1.47 * \text{Lat} + 97$$

$$X_{5x2} = \begin{bmatrix} 42 & 1 \\ 39 & 1 \\ 30 & 1 \\ 19 & 1 \\ 9 & 1 \end{bmatrix} \quad y_{5x1} = \begin{bmatrix} 33 \\ 38 \\ 62 \\ 67 \\ 83 \end{bmatrix}$$
$$(X^T X)^{-1} X^T y \rightarrow w_{2x1} = \begin{bmatrix} -1.47 \\ 97 \end{bmatrix}$$

# Let's Predict The Weather

The EECS 442  
Weather  
Channel

<u>City</u>	<u>Latitude</u>	<u>Temp</u>	<u>Temp*</u>	<u>Error</u>
Ann Arbor	42	33	35.3	2.3
Washington, DC	39	38	39.7	1.7
Austin, TX	30	62	52.9	10.9
Mexico City	19	67	69.1	2.1
Panama City	9	83	83.8	0.8

# Is This a Minimum Viable Product?

The EECS 442  
Weather  
Channel

The  
Weather  
Channel



*Pittsburgh:*

$$\text{Temp} = -1.47 * 40 + 97 = 38$$

*Actual Pittsburgh:*  
45



*Berkeley:*

$$\text{Temp} = -1.47 * 38 + 97 = 41$$

*Actual Berkeley:*  
53



*Sydney:*

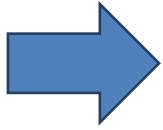
$$\text{Temp} = -1.47 * -33 + 97 = \mathbf{146}$$

*Actual Sydney:*  
74

Won't do so well in the Australian market...

# Where Can This Go Wrong?

# Where Can This Go Wrong?

Data			Model
<u>City</u>	<u>Latitude</u>	<u>Temp</u>	
Ann Arbor	42	33	
Washington, DC	39	38	$\text{Temp} = -1.6 * \bar{\text{Lat}} + 103$

**How well can we predict Ann Arbor and DC and why?**

# Always Need Separated Testing

Model might be fit data too precisely “*overfitting*”

Remember: #datapoints = #params = perfect fit

Model may only work under some conditions (e.g., trained on northern hemisphere).



Sydney:  
Temp =  $-1.47 \times -33 + 97 = 146$

# Training and Testing

Fit model parameters on **training** set;  
evaluate on *entirely unseen* **test** set.



“It’s tough to make predictions, especially about the future”  
-Yogi Berra

Nearly any model can predict data it’s seen. If your model  
can’t accurately interpret “unseen” data, it’s probably  
useless. We have no clue whether it has just memorized.

# Let's Improve Things

If one feature does ok, what about more features!?

<u>City Name</u>	<u>Latitude (deg)</u>	<u>Avg July High (F)</u>	<u>Avg Snowfall</u>	<u>Temp (F)</u>
Ann Arbor	42	83	58	33
Washington, DC	39	88	15	38
Austin, TX	30	95	0.6	62
Mexico City	19	74	0	67
Panama City	9	93	0	83

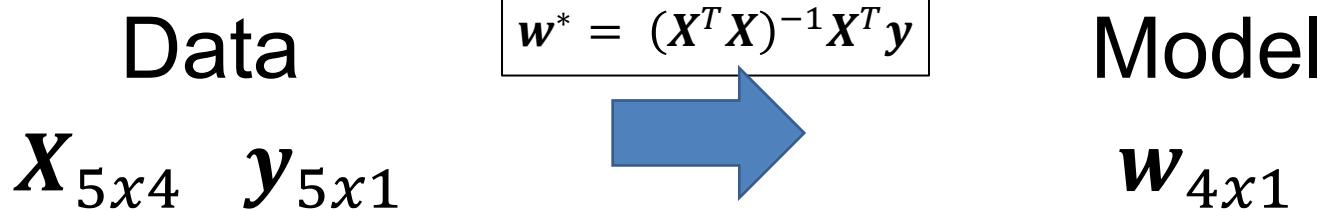
$$X_{5 \times 4}$$

3 features + a feature of 1s for intercept/bias

$$y_{5 \times 1}$$

# Let's Improve Things

All the math works out!



New EECS 442 Weather Rule:

$$w_1^* \text{latitude} + w_2^*(\text{avg July high}) + \\ w_3^*(\text{avg snowfall}) + w_4^*1$$

*In general called linear regression*

# Let's Improve Things More

If one feature does ok, what about **LOTS** of features!?

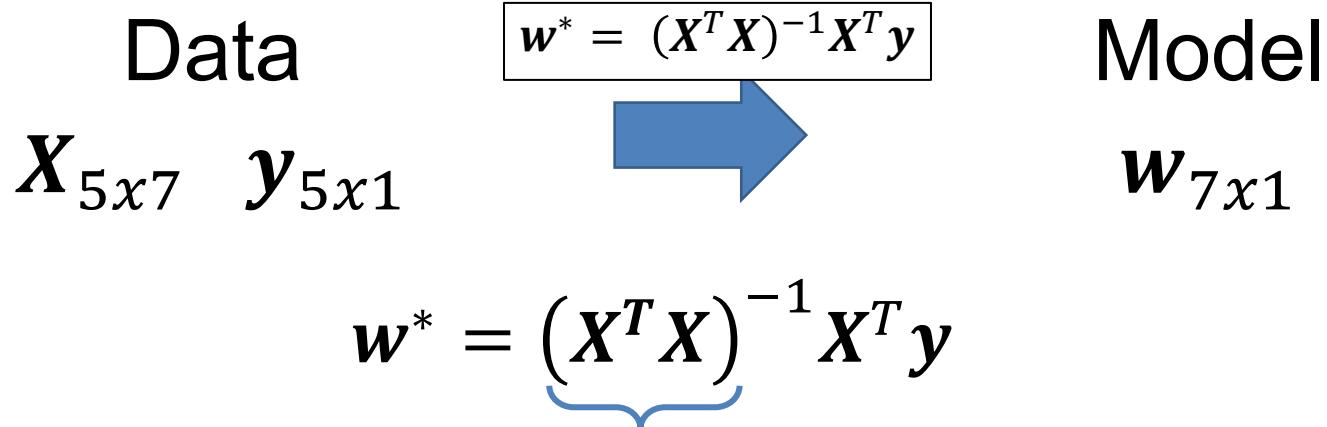
<u>City Name</u>	<u>Latitude (deg)</u>	<u>Avg July High (F)</u>	<u>Avg Snowfall</u>	<u>Day of Year</u>	<u>Elevation (ft)</u>	<u>% Letter M</u>	<u>Temp (F)</u>
Ann Arbor	42	83	58	45	840	100	33
Washington, DC	39	88	15	45	409	3	38
Austin, TX	30	95	0.6	45	489	2	62
Mexico City	19	74	0	45	7200	4	67
Panama City	9	93	0	45	7	1	83

$$X_{5 \times 7}$$

6 features + a feature  
of 1s for intercept/bias

$$y_{5 \times 1}$$

# Let's Improve Things More



$\mathbf{X}^T \mathbf{X}$  is a  $7 \times 7$  matrix but is **rank deficient** (rank 5) and has no inverse. *There are an infinite number of solutions.*

Have to express some preference for which of the infinite solutions we want.

Exercise for the mathematically-inclined folks: derive what the space of solutions looks like.

# The Fix – Regularized Least Squares

Add **regularization** to objective that prefers some solutions:

Before:  $\arg \min_w \|y - Xw\|_2^2 \longrightarrow \text{Loss}$

After:  $\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$

The diagram illustrates the components of the regularized least squares objective function. It shows the original loss term  $\|y - Xw\|_2^2$  in red, which is labeled 'Loss' with a red arrow below it. This term is combined with a regularization term  $\lambda \|w\|_2^2$  in orange, which is labeled 'Regularization' with an orange arrow below it. A blue double-headed arrow between the two terms is labeled 'Trade-off', indicating the balance between minimizing the loss and the regularization term.

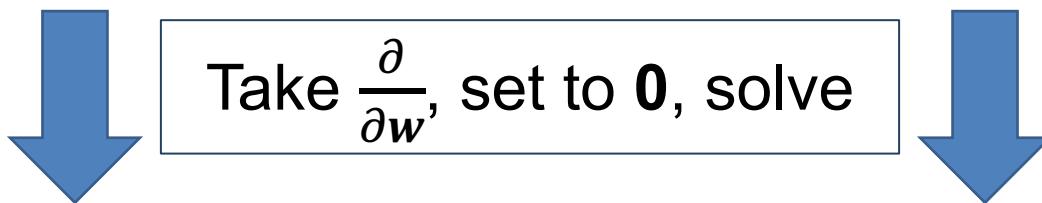
Want model “smaller”: pay a penalty for  $w$  with big norm

Intuitive Objective: accurate model (low loss) but not too complex (low regularization).  $\lambda$  controls how much of each.

# The Fix – Regularized Least Squares

Objective:  $\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$

The diagram illustrates the components of the objective function. A red arrow points from the term  $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$  to the word "Loss". A blue arrow points from the term  $\lambda \|\mathbf{w}\|_2^2$  to the word "Regularization". A blue arrow points from the term  $\lambda \|\mathbf{w}\|_2^2$  to the word "Trade-off".



$$\mathbf{w}^* = (\underbrace{\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}}_{\text{Matrix}})^{-1} \mathbf{X}^T \mathbf{y}$$

$\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is full-rank (and thus invertible) for  $\lambda > 0$

Called *lots of things*: regularized least-squares, Tikhonov regularization (after Andrey Tikhonov), ridge regression, Bayesian linear regression with a multivariate normal prior.

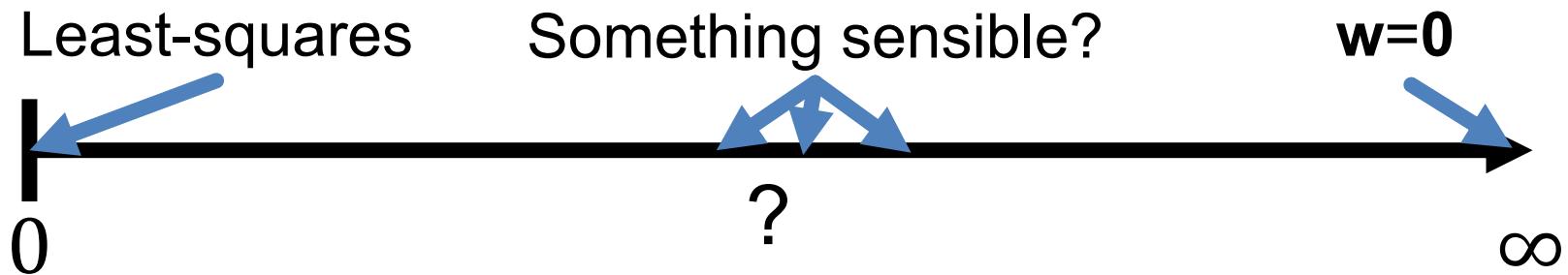
# The Fix – Regularized Least Squares

Objective:  $\arg \min_w \|y - Xw\|_2^2 + \lambda \|w\|_2^2$

Loss      Trade-off      Regularization

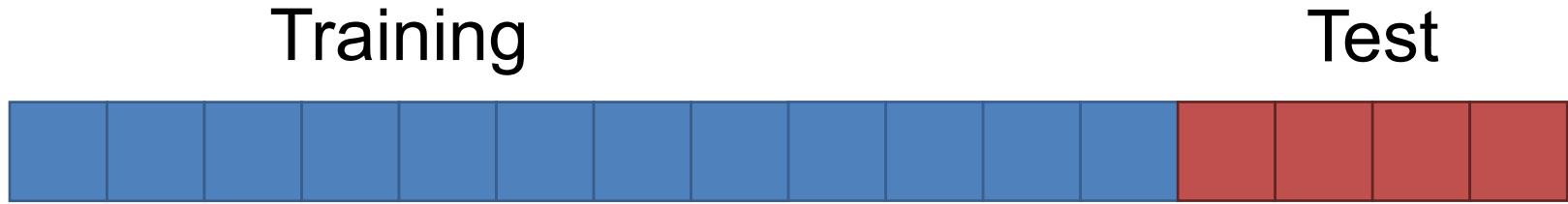
What happens (and why) if:

- $\lambda=0$
- $\lambda=\infty$



# Training and Testing

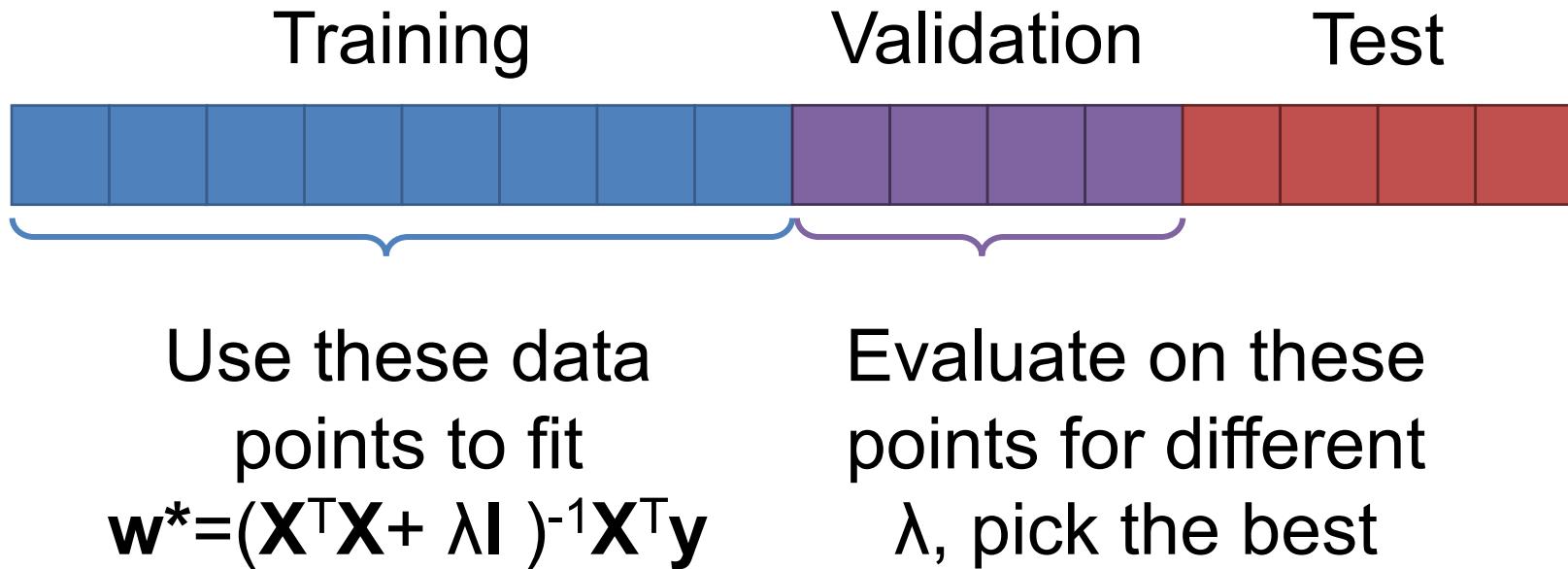
Fit model parameters on training set;  
evaluate on *entirely unseen* test set.



## How do we pick $\lambda$ ?

# Training and Testing

Fit model parameters on training set;  
find *hyperparameters* by testing on validation set;  
evaluate on *entirely unseen* test set.



# General Formula

Specific  
Formulation

$$\arg \min_{\mathbf{w}} \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|_2^2$$

General  
Recipe

$$\arg \min_{\mathbf{w}} R(\mathbf{w}) + \sum_{i=1}^n L(y_i, f(\mathbf{w}, \mathbf{x}_i))$$

Regularization

$$R(\mathbf{w}) = \lambda \|\mathbf{w}\|_2^2$$

Evaluation

$$f(\mathbf{w}, \mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

Loss

$$L(y_i, \hat{y}_i) = \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2^2$$

# Terminology

$$\arg \min_{\mathbf{w}} \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|_2^2$$

$x$  Inputs, features, Xs, data

$y$  Outputs, targets, labels, ys

$w$  Weights, weight vector, parameters, params

$\lambda$  Trade-off parameters, regularization coefficients

# Classification

Start with simplest example: binary classification



Cat or not cat?

Actually: a feature vector  
representing the image

What are some ways?

# Classification by Least-Squares

Treat as regression:  $x_i$  is image feature;  $y_i$  is 1 if it's a cat, 0 if it's not a cat. Minimize least-squares loss.

Training ( $\mathbf{x}_i, y_i$ ):

$$\arg \min_{\mathbf{w}} \sum_{i=1}^n \|\mathbf{w}^T \mathbf{x}_i - y_i\|^2$$

Inference ( $\mathbf{x}$ ):

$$\mathbf{w}^T \mathbf{x} > 0.5$$

Unprincipled in theory, but often effective in practice

Rifkin, Yeo, Poggio. *Regularized Least Squares Classification*  
[\(<http://cbcl.mit.edu/publications/ps/rpsc.pdf>\)](http://cbcl.mit.edu/publications/ps/rpsc.pdf). 2003

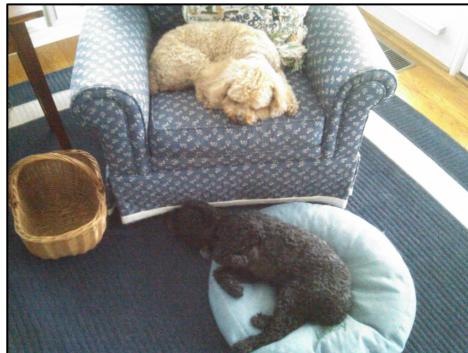
Redmon, Divvala, Girshick, Farhadi. *You Only Look Once: Unified, Real-Time Object Detection*.  
CVPR 2016.

# Easiest Form of Classification

Just **memorize** (as in a Python dictionary)  
Consider cat/dog/hippo classification.



If this:  
cat.



If this:  
dog.



If this:  
hippo.

# Easiest Form of Classification

Where does this go wrong?



Rule: if this,  
then cat



Hmmm. Not quite the  
same.

# Easiest Form of Classification

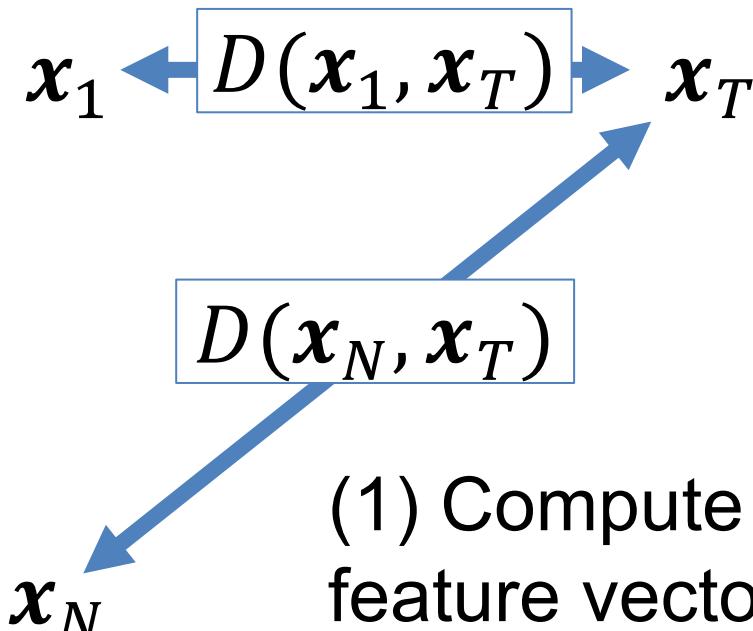
Known Images  
Labels



...



Test  
Image



- (1) Compute distance between feature vectors
- (2) find nearest
- (3) use label.

# Nearest Neighbor

“Algorithm”

Training ( $x_i, y_i$ ):

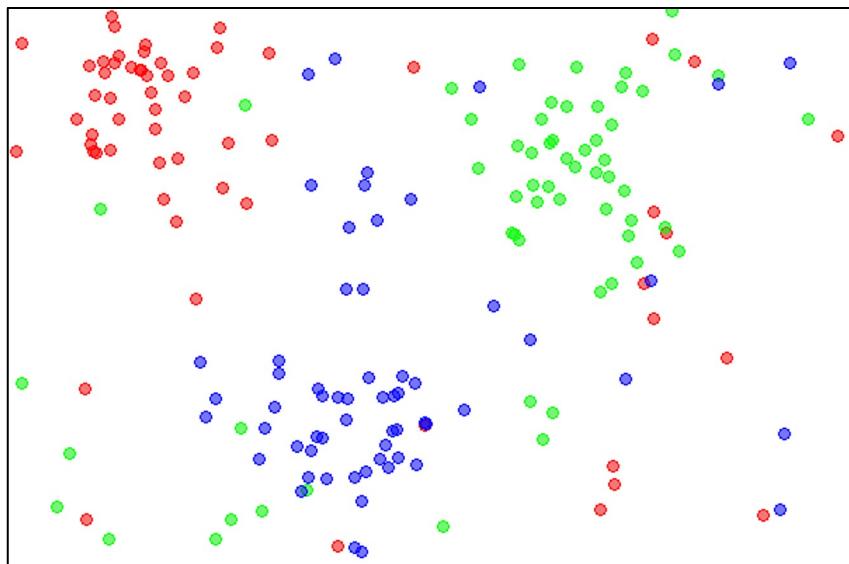
Memorize training set

Inference ( $x$ ):

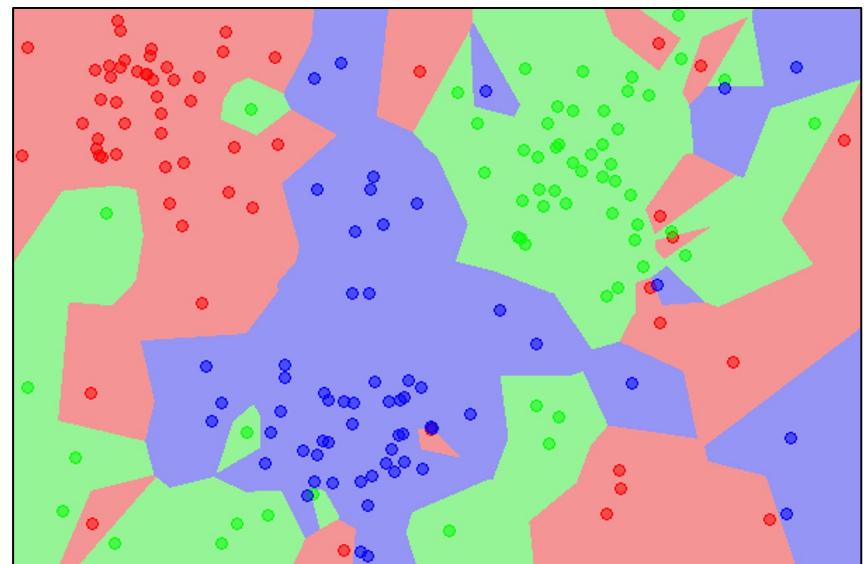
```
bestDist, prediction = Inf, None  
for i in range(N):  
    if dist( $x_i, x$ ) < bestDist:  
        bestDist = dist( $x_i, x$ )  
        prediction =  $y_i$ 
```

# Nearest Neighbor

2D Datapoints  
(colors = labels)



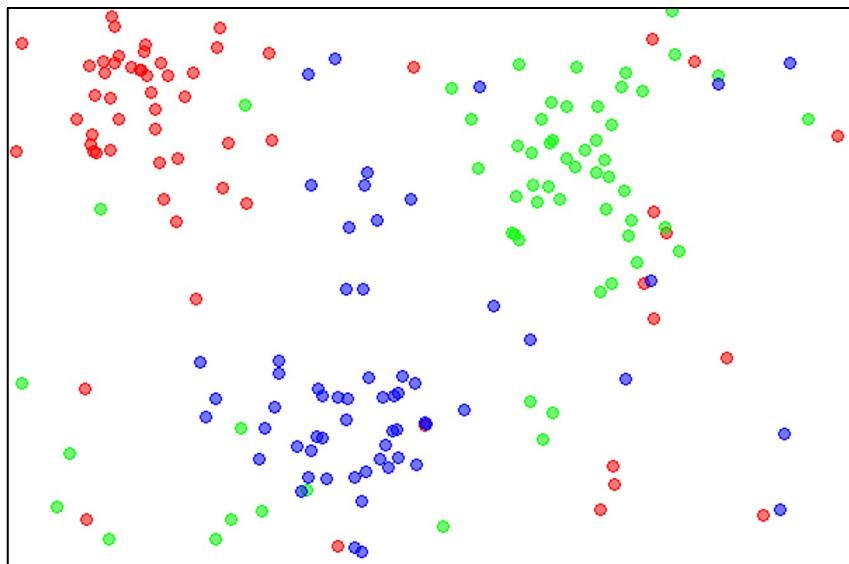
2D Predictions  
(colors = labels)



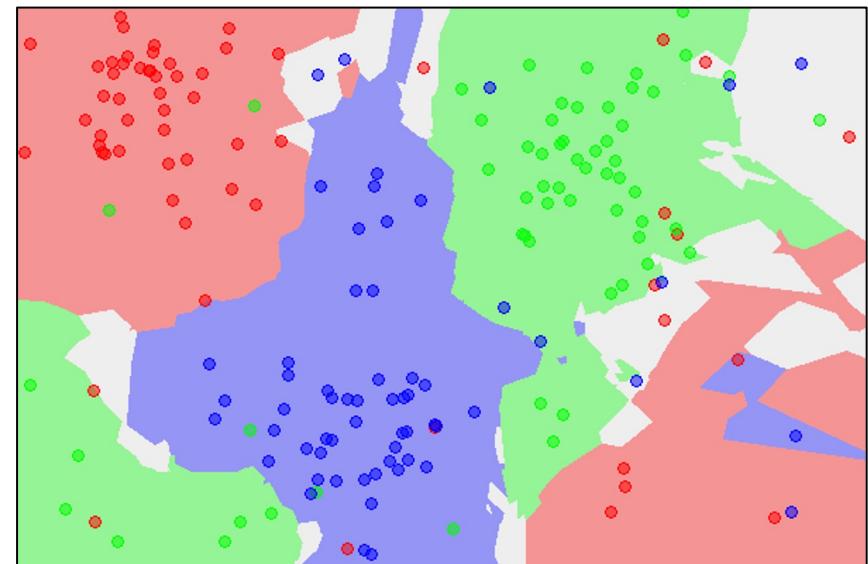
# K-Nearest Neighbors

Take top K-closest points, vote

2D Datapoints  
(colors = labels)

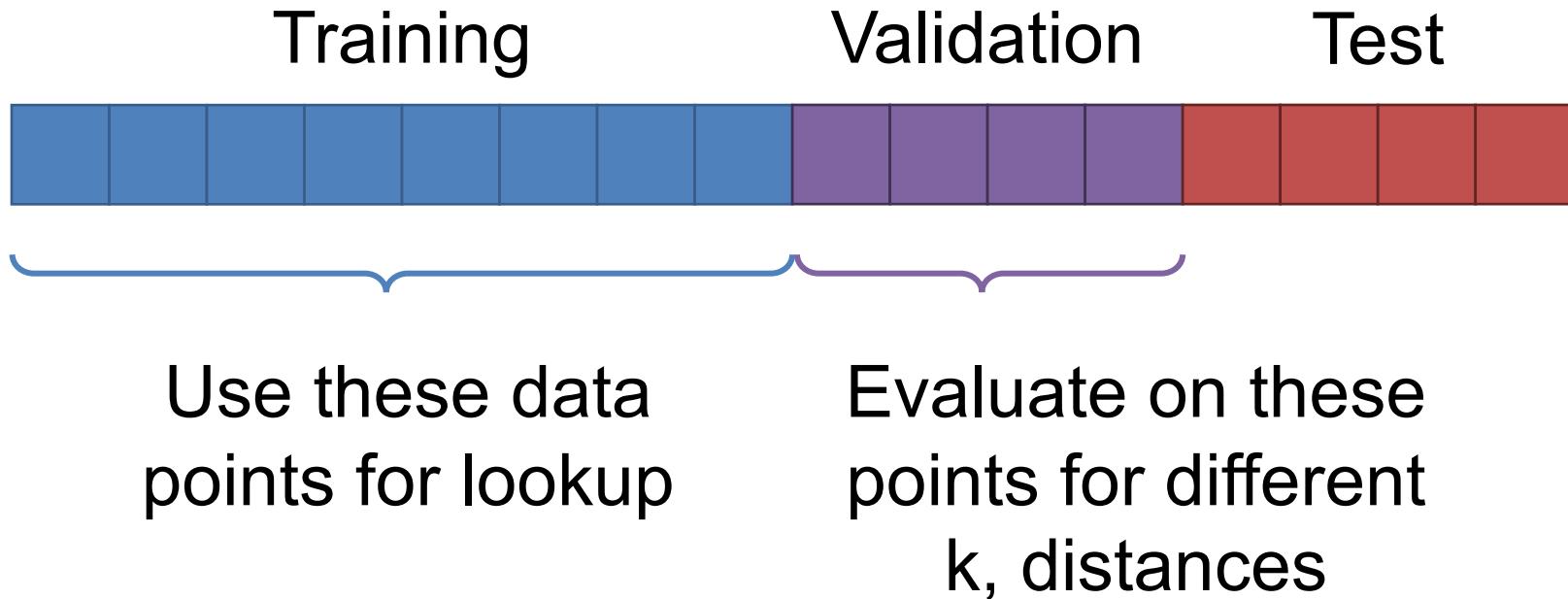


2D Predictions  
(colors = labels)



# K-Nearest Neighbors

What distance? What value for K?



# K-Nearest Neighbors

- No learning going on but sometimes effective
- Same algorithm for every task
- As number of datapoints  $\rightarrow \infty$ , error rate is guaranteed to be at most 2x worse than optimal you could do on data (error bound)
- Suffers from “Curse of dimensionality.” As feature dimension increases, the data becomes sparser, distance metric less useful.

# Linear Models

Example Setup: 3 classes



Model – one weight per class:  $w_0, w_1, w_2$

$w_0^T x$  big if cat

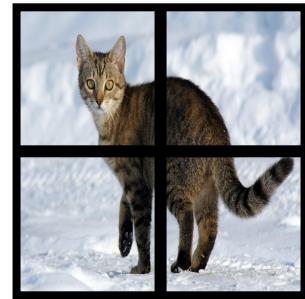
$w_1^T x$  big if dog

$w_2^T x$  big if hippo



Stack together:  $W_{3xF}$  where  $x$  is in  $R^F$

# Linear Models



Cat weight vector

0.2	-0.5	0.1	2.0	1.1
1.5	1.3	2.1	0.0	3.2
0.0	0.3	0.2	-0.3	-1.2

Dog weight vector

0.2	-0.5	0.1	2.0	1.1
1.5	1.3	2.1	0.0	3.2
0.0	0.3	0.2	-0.3	-1.2

Hippo weight vector

$W$

Weight matrix a collection of scoring functions, one per class

56
231
24
2
1

$x_i$

-96.8
437.9
61.95



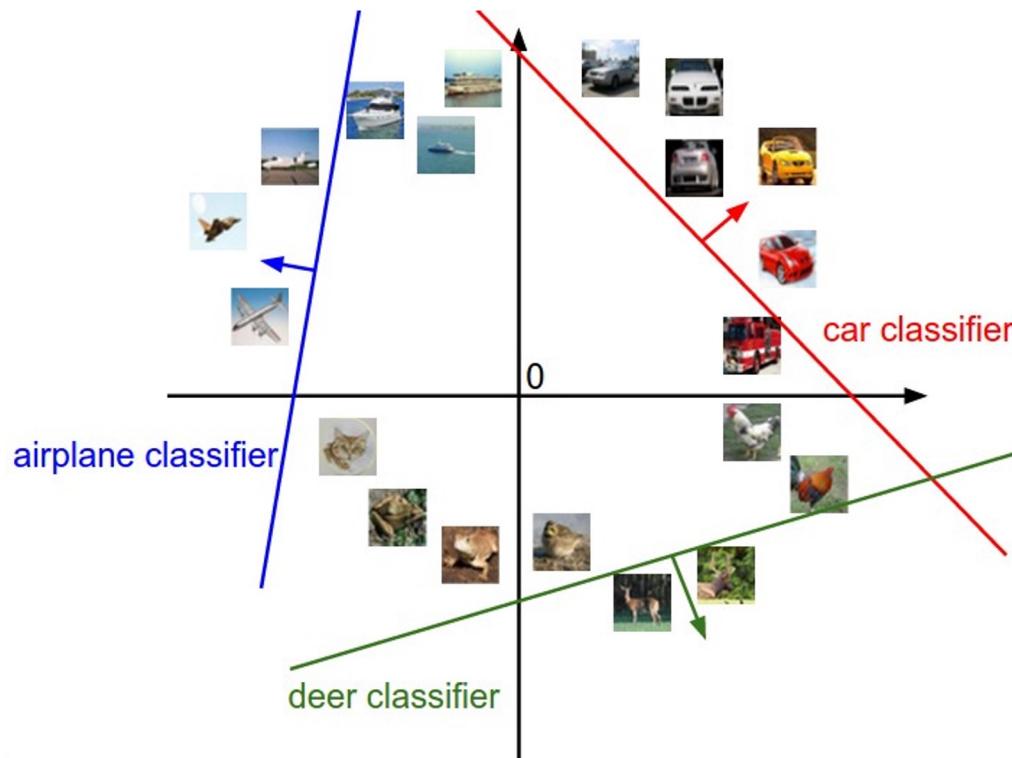
Cat score  
Dog score  
Hippo score

$Wx_i$

Prediction is vector where jth component is “score” for jth class.

# Geometric Intuition\*

What does a linear classifier look like\* in 2D?

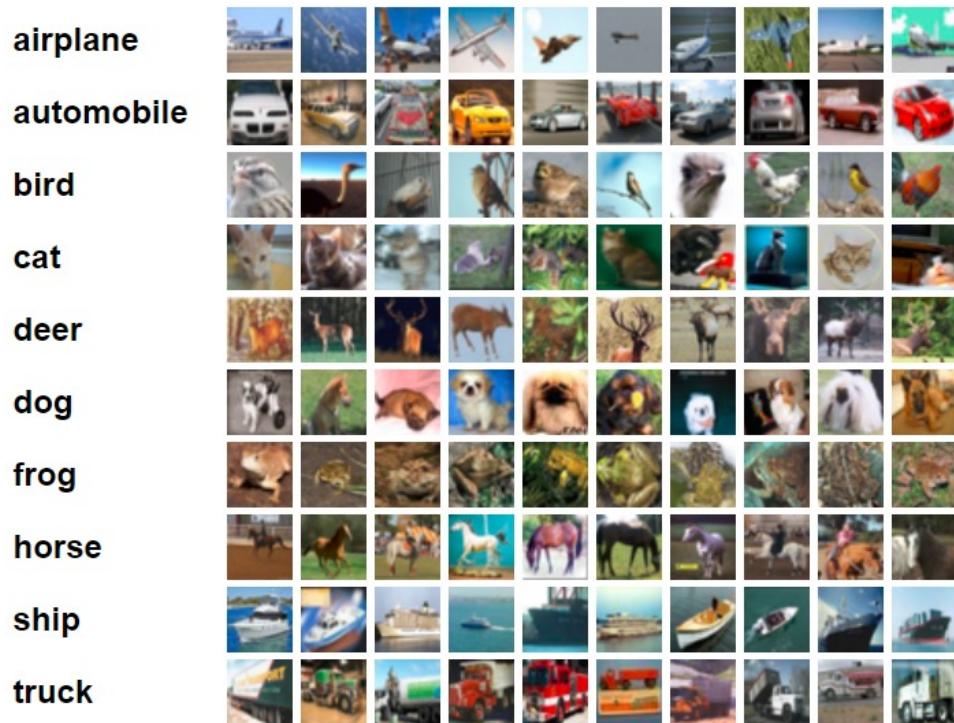


\*2D is good for vague intuitions, but ML typically deals with at least dozens if not *thousands* of dimensions. Your intuitions about space and geometry from living in 3D are **completely wrong** in high dimensions. Never trust people who show you 2D diagrams and write “Intuition” in the slide title. See: *On the Surprising Behavior of Distance Metrics in High Dimensional Space*. Charu, Hinneburg, Keim. ICDT 2001

Diagram credit: Karpathy & Fei-Fei. 12-point font mini-rant: me

# Visual Intuition

CIFAR 10:  
32x32x3 Images, 10 Classes



- Turn each image into feature by unrolling all pixels
- Fit 10 linear models

$$w, x_i \in \mathbb{R}^{32 \times 32 \times 3 + 1}$$

# Guess The Classifier

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then big values of  $x_i$  are indicative of the class.

## Deer or Plane?



# Guess The Classifier

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then big values of  $x_i$  are indicative of the class.

## Ship or Dog?



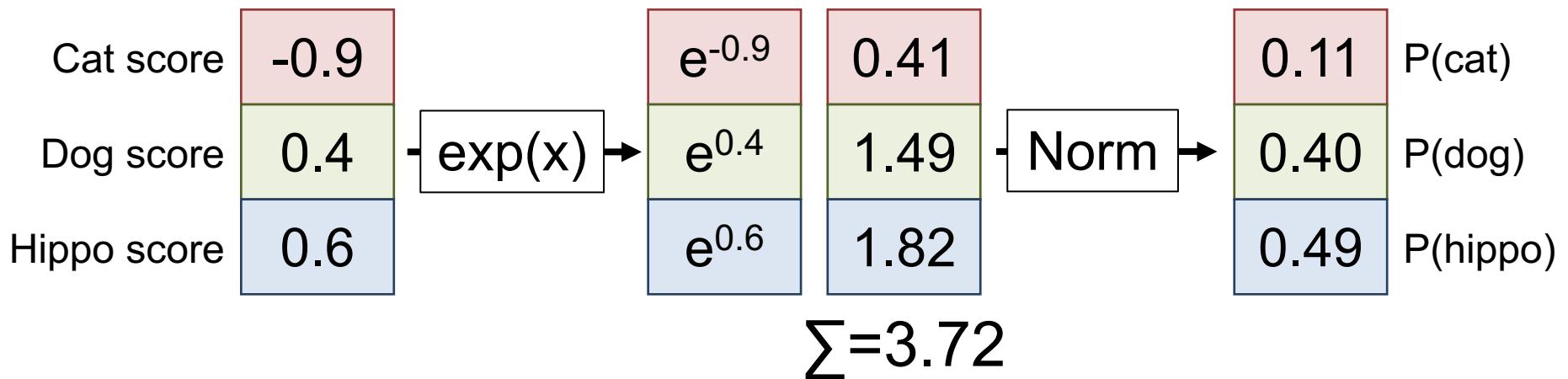
# Interpreting a Linear Classifier

Decision rule is  $\mathbf{w}^T \mathbf{x}$ . If  $w_i$  is big, then big values of  $x_i$  are indicative of the class.



# Preliminaries

Converting Scores to “Probability Distribution”



Generally  $P(\text{class } j)$ : 
$$\frac{\exp((Wx)_j)}{\sum_k \exp((Wx)_k)}$$

# Softmax

Inference ( $x$ ):  $\arg \max_k (Wx)_k$

(Take the class whose weight vector gives the highest score)

$$P(\text{class } j) = \frac{\exp((Wx)_j)}{\sum_k \exp((Wx)_k)}$$

**Why can we skip the exp/sum exp thing to make a decision?**

# Softmax

Inference ( $x$ ):  $\arg \max_k (\mathbf{W}x)_k$

(Take the class whose weight vector gives the highest score)

Training  $(\mathbf{x}_i, y_i)$ :

$$\arg \min_{\mathbf{W}} \lambda \|\mathbf{W}\|_2^2 + \sum_i -\log \left( \frac{\exp((\mathbf{W}x)_{y_i})}{\sum_k \exp((\mathbf{W}x)_k)} \right)$$

Regularization  $\lambda \|\mathbf{W}\|_2^2$  (Orange arrow)

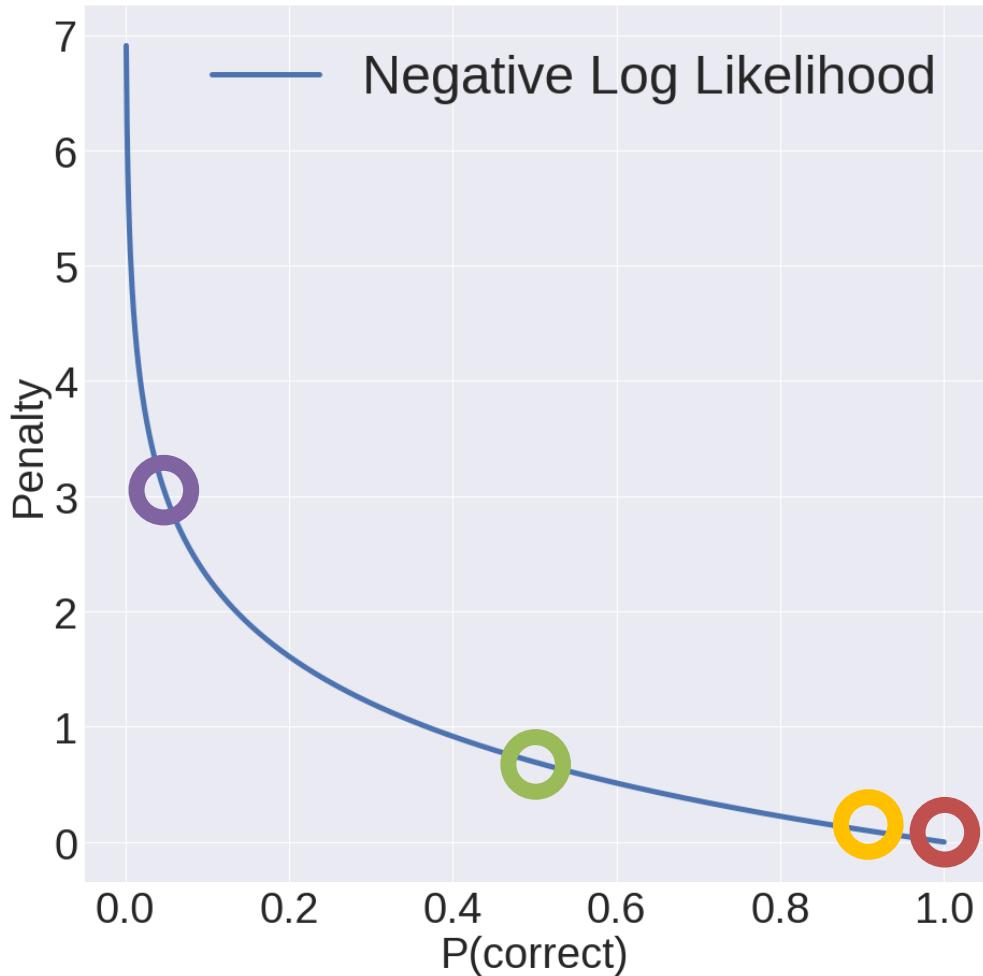
Over all data points  $\sum_i$  (Teal arrow)

Pay penalty for negative log-likelihood of correct class (Purple bracket)

P(correct class) (Red arrow)

No Closed-Form Solutions

# Negative Log Softmax



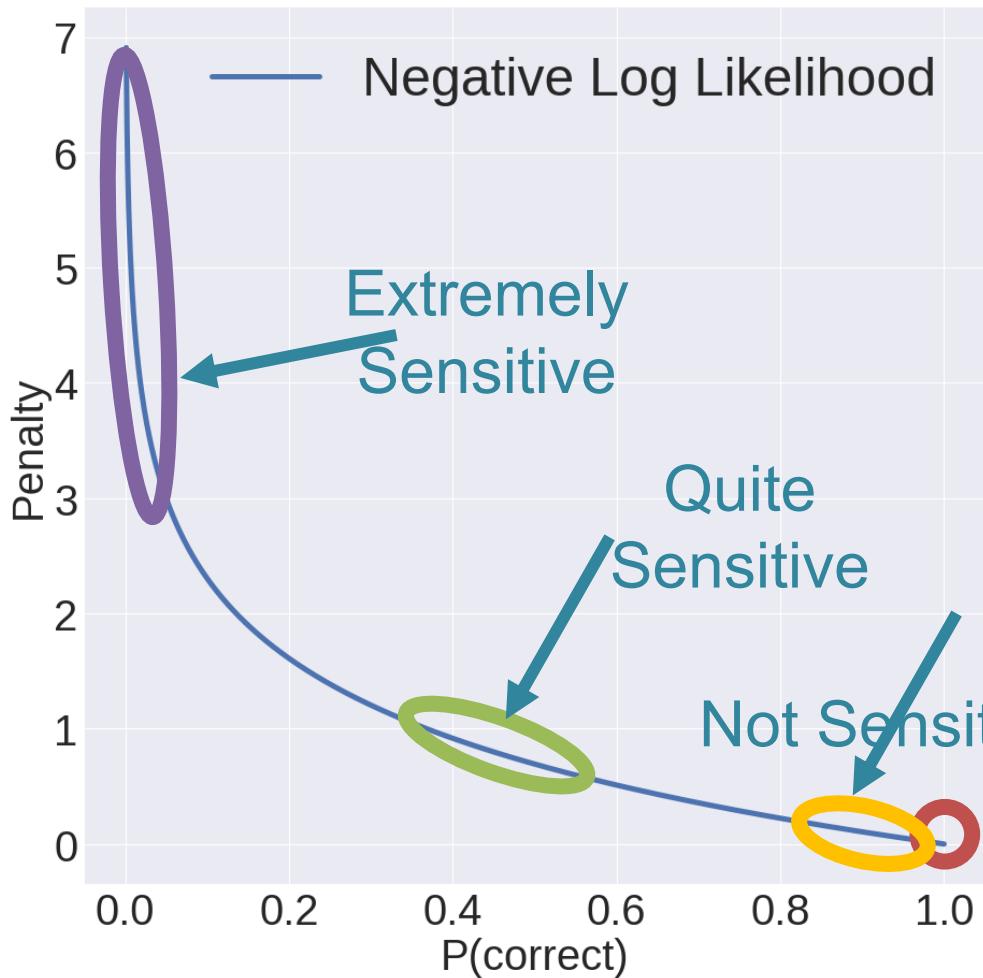
$P(\text{correct}) = 0.05:$   
3.0 penalty

$P(\text{correct}) = 0.5:$   
0.7 penalty

$P(\text{correct}) = 0.9:$   
0.1 penalty

$P(\text{correct}) = 1:$   
No penalty!

# Negative Log Softmax



$P(\text{correct}) = 0.05$ :  
3.0 penalty

$P(\text{correct}) = 0.5$ :  
0.7 penalty

$P(\text{correct}) = 0.9$ :  
0.1 penalty

$P(\text{correct}) = 1$ :  
No penalty!

# Next Class

- How do we optimize more complex stuff?
- A bit more ML



# Objective 1: Multiclass SVM

Inference ( $x$ ):  $\arg \max_k (\mathbf{W}x)_k$

(Take the class whose weight vector gives the highest score)

# Objective 1: Multiclass SVM

Inference ( $x$ ):  $\arg \max_k (\mathbf{W}x)_k$

(Take the class whose weight vector gives the highest score)

Training ( $\mathbf{x}_i, y_i$ ):

$$\arg \min_{\mathbf{W}} \lambda \|\mathbf{W}\|_2^2 + \sum_i^n \sum_{j \neq y_i} \max(0, (\mathbf{W}x_i)_j - (\mathbf{W}x_i)_{y_i} + m)$$

Regularization

Over all data points

For every class j that's NOT the correct one ( $y_i$ )

Pay no penalty if prediction for class  $y_i$  is bigger than  $j$  by  $m$  ("margin"). Otherwise, pay proportional to the score of the wrong class.

# Objective: Multiclass SVM

How on earth do we optimize:

$$\arg \min_{\mathbf{W}} \lambda \|\mathbf{W}\|_2^2 + \sum_i^n \sum_{j \neq y_i} \max(0, (\mathbf{W}\mathbf{x}_i)_j - (\mathbf{W}\mathbf{x}_i)_{y_i} + m)$$

Hold that thought!

# Softmax Temperature

$$P(\text{class } j) = \frac{\exp(x_j/T)}{\sum_k \exp((Wx)_k/T)}$$

When  $T \ll 1$ ,  $\text{Softmax}(x, T) \Rightarrow \arg \max_k x_k$

Softmax is essentially “soft argmax”